

A Review on Congestion Control Approaches for Real-Time Streaming Application in the Internet

Abhay Kumar

Department of CSE, J B Institute of
Engineering and Technology,
Hyderabad, India
abhay25.cse@gmail.com

P. V. S. Srinivas

Sreenidhi Institute of Technology and
Science, Hyderabad, India.
pvssrinivas23@gmail.com

Dr. A Govardhan

Department of CSE, Jawaharlal
Nehru Technological University
Hyderabad, India
govardhan_cse@yahoo.co.in

Abstract— In the support of congestion control over the Internet in providing the assurance of the equality between much diverse traffic is a difficult function. The advent of streaming media has offered users with low-latency media content, with higher congestion on the Internet due to stringent bandwidth and latency requirements. Therefore, it is more and more important to resolve the difficulties of increased packet deliver fail reasoned because of congestion and better quality of service for streaming media. In this paper, we propose a review on the congestion control approaches (CCA) for the real-time streaming applications on the Internet. The role of TCP in network congestion control and the characteristics of the original real-time streaming media are discussed. After that, we discuss issues in the media stream and real-time congestion control. The survey will support the understanding of the current congestion mechanism and continue to enhance the expansion of real-time streaming application services.

Keywords- Internet, Congestion Control Approaches, Real-Time Streaming, Multimedia

I. INTRODUCTION

The Internet presents end-to-end finest-attempt for data packet services through utilizing IP protocol exclusive of any unambiguous assurance to the quality of service or deliverance assurance. The presentation and steadiness of the Internet depend on the congestion control approaches (CCA) applied in the "transport layer" protocol. The principal function of transport layer is to utilize internet through "Transmission Control Protocol (TCP)", which consists of a congestion window algorithm that runs on the terminal system. As long as no packet is lost, TCP progressively enhances the transfer rate of the traffic source. This allows traffic sources to establish to what extent the bandwidth is accessible in the network without congestion. Packet deliver fail happens while the network is congested. In response, TCP reduced the sending rate to control congestion. Therefore, TCP adds a multiplication-reduction (AIMD) policy [1] by using addition to cooperatively adjust the sending rate from the suggestion of network congestion in the structure of discarding packets. This approaches permits traffic sources to distinguish congestion and "back off" to accomplish a central delivery rate which will be equal to the congestion point's competence. Therefore, a TCP flow is referred to as a "response" [2] congestion indication (eg, a dropped packet) from the network because congestion control reasons the TCP transmitter to reverse off when congestion is identified.

In these day's internet-based real-time applications such as "VoIP", "video conferencing", and "online games" predominantly use "RTP over UDP" or "UDP" to transmit data. As these protocols do not respond to the congestion measures, the use of these protocols is gaining in popularity, compromising Internet stability [1], [7]. Therefore, in order to construct the real-time application extensively used, it is estimated that a general CCA suitable for real-time multimedia will be deployed. The Internet tends to support the maximum transfer rate of traditionally supported applications, whereas media streaming applications require smoother transfer rates and less jitter. Therefore, in order to securely deploy streaming media applications over the Internet, new congestion control algorithms need to be developed to allow fair interactions with other TCP-oriented applications and maintain the steadiness of the Internet.

In the case of network users competing for scarce network bandwidth, the deliver fail rate is very high during severe congestion in data networks [6]. Internet surveys show a growing demand for bandwidth-intensive applications, resulting in an increase in the delivery failure rate across the entire characteristics of the Internet. It should take appropriate approaches to control network congestion, or the network may be in sustained overload, which may lead to the collapse of the network. Implicit CCA employed in the main transport protocols, TCP used by HTTP and FTP, help improve the robustness of the Internet [1], [5], [9]. Delay-sensitive media applications such as media streams and Webcasts that do not use TCP have disrupted this friendly best-effort network by not responding to network congestion [3], [4]. Many algorithms have been recommending to complement TCP's friendliness for uninterrupted media applications [8], [9], [11], [13], [14].

However, the Internet needs to provide some outline of feedback to data traffic originating from the congested links so that it can regulate its transmission rate depending on the accessible bandwidth, effectively managing end-to-end congestion control [16]. Feedback on congestion can be implicit or explicit. In the case of implicit response, the network's transport layer protocol attempts to maintain high throughput by approximation "service time", "end-to-end delay", and "packet deliver fail". The TCP protocol widely used by the Internet [7], [10] implicitly feeds back lost packets over time and repeatedly. Terminal nodes usually deploy explicit feedback. However, relying on end nodes for implicit

or explicit feedback is not enough to achieve a high throughput of the Internet.

The purpose of this survey is to design effective data routing mechanism to detect congestion and reduce packet deliver fail rate efficiently to advance the competence of Internet throughput, to decrease the packet deliver fail rate and to diminish the network bandwidth requirements. The router's efficient routing mechanism can support the reduction of end-to-end congestion control, which is a major concern for Internet traffic and manages high-bandwidth traffic during congestion. It helps application developers and protocol designers provide the best congestion control for Internet traffic.

The following paper is organized to discuss the TCP in internet congestion control in section-2, Real-time streaming on the Internet in section3, Congestion Control in media streaming mechanism in section-4, investigation of related works in section-5 and the conclusion of the paper in section-6.

II. TCP IN STREAMING CONGESTION CONTROL

During high traffic rates, most TCP methods for congestion control cause a reduction in rate. This results in a huge change in send rate, which is identified by the difference in video quality or a large delay in data buffering before being played, creating a large waiting time for the user. TCP ensures very strict reliability and ordering semantics at the expense of end-to-end delay, which is important for legacy applications such as file transfers but when the data is being played to the user in real time, the receiver is Useless. Even responses to dropped packets are retransmitted, resulting in higher network load and significantly lower effective throughput.

TCP is not suitable for emerging applications that stream multimedia content. Media streaming applications are both a data transfer rate and a data transfer rate change. When a media streaming application presents interactive media data to a user in real time, changes in the reception rate at the receiver are visible to the user. These changes are handled by buffering at the receiver, but waiting for the desired buffering results delays the response, uncomfortable viewing, and poor interactivity. Because TCP is purely window-based, it can cause data bursts that get worse due to acknowledged transmissions. This bursty nature makes it very difficult and inefficient to preserve the relative timing of the various data frames received at the receiver [14], [20], [21].

Therefore, TCP with implicit end-to-end congestion control using convective multimedia content is not preferred and a modified protocol with explicit congestion control and the following issues must be used:

- Real-time transmission of media data over the network necessitates high network output because it is simplified to make up for missing data than to make up for the greater delay in reception media data. This is not the case with regular data such as files, where the complete arrival of data at its destination is necessary, so protocols for these static data are not suitable for streaming media.

- In today's computer networks, more and more multimedia applications cannot utilize "Transmission Control Protocol (TCP)" congestion control as a result of "stringent latency" and "jitter" necessities. As the level of congestion enlarges, the excellence degradation of these multimedia functions eventually achieves the level at which users cannot receive content.

The conventional deliver fail-based TCP is not appropriate for real-time traffic because its congestion control continually detects the network's accessible bandwidth, establishing the intervallic rotations throughout which network lineups are foremost packed and then it flow down. These queue fluctuations cause random delay components that vary over time, which increases the circulation time and constructs "delay-sensitive communications" problem. There are two complementary ways to solve this problem: "End-to-End Delay", "Control Over Endpoints", and "Active Queue Management (AQM)" to solve problems in routers [12], [26].

In the following, we provide comments on the end-to-end congestion control algorithm proposed by the related work clustering based on the measures used to derive the congestion and the AQM to control traffic jam queuing setbacks in the network.

A. The exploit of RTT to Infer Congestion

In the field of TCP congestion control, a method of reducing the queuing delay is first designed. Therefore, many real-time business algorithms are based on the past literature studies. The initial congestion control algorithm was particularly applied in Jain's groundbreaking work to include end-to-end delay, dating back to 1989 [35]. Since then, numerous delay-based TCP congestion control modifications comprise are proposed, for example, "TCP FAST" [27] and "TCP Vegas" [36], that utilize "RTT measurements" to conclude congestion. It has been shown that when the RTT is used as a congestion metric, low channel utilization can be achieved when there is reverse traffic or competition with deliver fail-based traffic [19]. It is significant to talk about that in the video conferencing environment, the issue of back traffic is essential because the video stream is sent in both directions.

B. The utilize of Delay-Gradient to Infer Congestion

The initiative of using "RTT gradients" to understand congestion has recently been adopted to prevail over the "latecomer effect" described above. Few instances are "Verus" [23] and "CDG" [28]. The "Verus" [23] is designed specifically for cellular networks and changing the burst link capacity makes congestion control design challenges. The "CDG" [26] aims to offer reasonable subsistence of "loss-based traffic" and "low end-to-end latency". Recently, accurate delay gradient measurements have been demonstrated in data center networks using "NIC hardware time stamping" [29].

C. AQM Algorithms to Reduce the Queuing Delays

The queuing delay can also be minimized by adjusting the network buffer size appropriately [31], [32], [33] or using the AQM algorithm, which controls router buffers through

"reducing packets" or "marking" them as ECNs are employed [34]. Although numerous AQM algorithms have been proposed in the earlier period, their acceptance is blocked due to two major problems: (i) they intend to control the standard queue length as a substitute of the lineup delay, and (ii) the ad-hoc limitations must be configured. These problems, as well as the phenomenon of buffer expansion [30], facilitate the revision of the innovative AQM algorithm, for instance, "CoDel" [26] and "PIE" [24], which do not necessitate constraint modification and unambiguously control the queuing delay rather than the queue length.

III. REAL-TIME STREAMING ON INTERNET

The widespread features of media streaming applications consist of the necessitate for "high bandwidth", "smooth data flow", and "low predictability of end-to-end latency" and "latency differences" [3], [9], [10], [22]. In disparity, the Internet is a most excellent-effort network that does not provide superiority of service assurances, large variations in "network bandwidth" and "host processing speed", the extensive distribution of resources, and high workloads. So, in fact, Internet-based applications have undergone major changes in the available bandwidth, latency, and latency differences. Real-time transmission of media data over the network necessitates high network outcomes because it is simple to make up for lost data than to make up for the greater delay in reception of media data. This is not the case with regular data such as files, where the complete arrival of data at its destination is necessary, so protocols for these static data are not suitable for streaming media. Like "TCP", "UDP" it is also common transport layer protocol that establishes more proprietary protocols over this protocol.

The Internet paradigm for transmitting real-time data such as "audio" and "video" is "Real-time Transport Protocol" (RTP). It is integrated into the application-layer "Real-Time Streaming Protocol (RTSP)" that enables the deliverance of multimedia over the Internet. RTSP supports interoperability between media clients and media servers from different vendors. RTSP uses streaming media to decompose media data into multiple packets and adjusts supported over the accessible bandwidth between the client and the server.

In a typical network, congestion occurs when packets carried by the link exceed the bandwidth competence of the network. It can be described as two categories of congestion difficulties that may occur in the network [16], [17].

- In a "single bottleneck problem", numerous transmitters send their data packets to a destination over a router with an imperfect yield of bandwidth. This congestion knows as "dumbbell-shaped congestion problem". Realizing that congestion is primarily reasoned through the "one-way flow" of numerous transmitters.
- In a common "bottleneck problem", numerous transmitters send packets to several recipients over a series of routers. In the transmitter link, its bandwidth cannot meet the data transmission load requirements. It shows that congestion is a butterfly crowded dilemma. In the mutual "bottleneck problem", various recipients can

correspond with numerous transmitters. This can cause the connection to stream in both directions at the same time.

Previous evaluations [18], [19] argue that single bottlenecks, congestion will lead to higher packet deliver fail rates and share bottlenecks. On the other hand, when the traffic load is constant for a moment, congestion will result in low network throughput for bottleneck sharing issues. However, when the traffic load increases, it will also outcome in a higher packet delivery fail rate. The two congestion problems have a common effect, that buffer overflow. As a result, most CCA focus on resolve this "buffer overflow problem".

The "Real-time Control Protocol (RTCP)" is the object of "RTP", providing management services for streaming applications. The main utility of "RTCP" is to present reaction to the excellence of distribution data. The "RTCP" is used for session control, QoS reporting, and media synchronization. Adaptive multimedia transmitters and receivers are designed so that the protocol for streaming media is simulated. The multimedia receiver sends feedback to the transmitter on network congestion, and the transmitter adjusts its transmission rate accordingly.

IV. CONGESTION CONTROL IN MEDIA STREAMING MECHANISM

This section aims to discuss different congestion prediction and routing models to ensure the highest quality of media delivered under a given network condition. It also emphasizes stochastic transmission schemes as part of the model to reduce jitter and burst loss in media transmissions.

Congestion control in media streaming protocols is designed with rate control to ensure that traffic does not exceed congestion-sensitive TCP traffic, which forms a major part of Internet traffic. A TCP-friendly process has been devised that does not use extra bandwidth but responds to consistent TCP connections based on congestion notifications. Therefore, the object of the invention is to provide a rate adaptation mechanism that does not cause TCP traffic in the background to get in trouble and adjust the transmission rate accordingly to achieve "end-to-end congestion control". The "End-to-end congestion control" depends on the acknowledgment from the receiver so that the sender changes its transmission rate. The congestion control mechanism ensures that a TCP connection using AIMD gets its fair bandwidth allocation when congested.

Modifying the router mechanism includes congestion control at the network layer. If the incoming data rate is higher than the outgoing data rate, the router will not be able to accommodate newly arrived packets because there is no buffer space. In this case, the router must decide to drop the incoming packet. Many strategies are used to make this decision. The simplest and most widely deployed is the tailing algorithm, where each packet arriving at the router is queued until the buffer is full. If there is no space unique, new packets arriving at the router are discarded. Although the tailing algorithm is simple to implement, it shows serious interactions with TCP's congestion control mechanism, resulting in poor performance. Trailing might lead to global synchronization and locking.

Another strategy used in routers is "Active Queue Management", which discards packets before the buffer is full, allowing the source to react to initial congestion. RED is a widely deployed active queue management technology that randomly drops packets depend over the standard queue dimension, where the "queue length" is between the least and highest thresholds. RED is not designed to operate on any particular protocol but rather treats the protocol as a sign of congestion to better enforce the protocol.

V. INVESTIGATION OF RELATED WORKS

The limitations discussed above prompt us to propose ways to improve multimedia streaming by achieving better throughput and lower network load through effectively controlling traffic rates, congestion control and reducing retransmissions.

The "Dynamic adaptation" is a dominant approach, although it needs innovative traffic and CCA to accurately discover and properly utilize the obtainable network bandwidth. The CCA should dynamically establish the share of network bandwidth that adaptive applications be able to utilize moderately under competitive traffic conditions [37], [38]. If these approaches are not sensitive enough to challenging for traffic, probability high multimedia data rates can reason for severe network congestion.

Based on the identified TCP congestion control limits [3], [5], [20], this research work on Multimedia Streams will serve as an enhancement to contribute to the following goals:

- An end-to-end strategy is a classic approach to congestion control. In this model, the TCP sender must detect congestion and take action. Thus, all end-to-end congestion approaches can only rely on implicit congested RTT signals, i.e. packet delivers fail and delay variation. However, in the presence of delivering fail-based algorithms, there are many problems with accurate RTT estimation and lack of link utilization. It will contribute to new improvements in congestion prediction and control routing approaches based on end-to-end CCA. This improves multimedia streaming end-to-end latency and minimizes network overhead.
- Another approach to routing assisted congestion control algorithms will help to overcome the problem of tail-queuing. In the drop tail queue, packets are dropped on the bottleneck router after a queue buffer overflow occurs. In the RED management queue, packets will be discarded earlier based on the output of the discard probability algorithm. If queue usage is low, there is almost no packet deliver fail. Conversely, if the queue begins to populate, the drop probability increases proportionately.
- Finally, it will provide a novel TCP fairness queue mechanism for maintaining fair queue management for each flow that can control flow, and for handling flows that can saturate or delay flow to unacceptable flows. The "Fair queuing mechanisms" [12], [25] distribute the bandwidth of the line fairly and provide relatively small

queue delays for short communications to increase throughput.

G. Carlucci et. al. [1] proposed a new congestion control algorithm for RTC depended on the foremost thought of estimation - using "Kalman filter" in the end-to-end unidirectional stoppage dissimilarity understanding by the packet from the transmitter to the destination. It compares the approximation with a "dynamic threshold" and compel the forceful of the controller located at the receiver, which is intended to maintain a low queuing delay, and a deliver fail-based controller at the transmitter when the loss is found It works. It is adopted by "Google Chrome" for the congestion control. A large number of experimental evaluations show that this algorithm includes queuing delay while providing fairness of the intro and inter-protocols and full link utilization.

S. D'Aronco et al. [2] proposed in favor of "delay-constrained communication" on the most excellent "packet-switched networks" in a new congestion control algorithm. The algorithm maintains a restricted queuing delay when challenging with erstwhile delay-depend streams, avoiding malnourishment while challenging in favor of delivering fail-based streams. It uses the distinguished value-dependent allocated method as congestion control, however: (1) introduces a recent "non-linear mapping" between "empirical delay" and "value functions", (2) combines delay and loss in the sequence of packets based on packet arrival interval Single-price measurement period.

A. Biernacki [6] studied the traffic commencing "120 client-server pairs" in simulated lab surroundings and multiplexed against a particular network link. It shows that the arrangement of traffic is different from that of the "first generation" and "second generation HTTP video systems" and is not similar to the general Internet traffic construction. The traffic volume acquires shows negative correlation and anti-determination, and its distribution purpose is inclined to the right. In addition, it indicates that the traffic produced by users using the same or similar playout policies is positively related and synchronized (clustering), while the traffic from dissimilar play-out policies demonstrates negative or no correlation.

G. Tian et al. [9] formally studied the trade-off between responsiveness and smoothness in "HTTP Adaptive Streaming (HAS)" by analyzing and experimenting. This shows that client-area buffering of video moment is a superior response indication to conduct video alteration. Then proposed a new video rate control algorithm for stability the video rate smoothing and high bandwidth consumption requirements. It shows that a miniature video rate edge is able to result in a dramatic increase in the smoothness of video rates and buffer sizes. It also proposes HAS invents that exertion with multiple servers and wireless associations. This proves that our HAS design is well-organized and vigorous in a practical network impression.

Most of the existing work [1], [2], [4], [5], [7] usually adopts the method of random early detection to adjust the source rate of the video sender in point of network congestion feedback and maintain the TCP friendliness when the network is congested The transmission. To overcome this problem, most of these models enable the interaction between Sender-

adaptive transmission and router-first packet filtering to achieve low loss, high-quality multimedia delivery over the Internet.

In order to diminish packet deliver fail reasoned through "bandwidth searching", "video congestion control algorithms" can utilize self-assured fault modification to discard dropped packets. Regrettably, if redundancy is added to the dominant part of the contention flow, additional overhead may cause an increase in packet delivery fail in the network, resulting in additional destruction. However, in this article, it has been established through analysis and simulation that the negative impact of balancing additional overhead can be heightened by means of a paradigm of "DiffServ" discard method and handover inferior priority to "FEC packets". The analysis also shows that the benefits of FEC are maximized when the router buffer is small. This is well related to the latency and jitter requirements of multimedia streams. Transferring video above the Internet is a significant part of numerous multimedia applications. Currently, the Internet lacks QoS support. The heterogeneity of network and terminal systems creates several confronts for the invent of video transmission systems.

In order to offer readers an apparent understanding of this intend gap, we review the improvements and drawbacks of the following methods and solutions.

- 1) Congestion control: There are three CCA: "rate control", "rate adaptive video coding", and "rate shaping". The "Rate control methods" fall into three kinds: "source-based", "receiver-based", and "hybrid". The rate control scheme is able to go after a "model-based approach" or a "Probe-based approach". The "Source-based rate control" is mainly for unicast and be able to be based on a "model-based approach" or a "Probe-based approach". If functional to "multicast", "source-based rate control" be able to simply go after "Probe-based approaches". The "Source-based rate control" requires a further element to implement the rate on the video stream. This element can be "rate-adaptive video coding" or "rate shaping". For illustrations unite "source-based rate control" and "rate-adaptive video coding" be able to be establish in [51], [63]. For instances of "source-based rate control" and "rate shaping" take account in [25]. It is proposed to solve the problem of heterogeneity in multicast video based on receiver and hybrid rate control.

The benefit of "receiver-based control" more than "sender-based control" is that the trouble of alteration moves from the transmitter to the receiver, thereby increasing the flexibility and scalability of the service. The "Receiver-based rate control" may go after a "model-based approach" or a "Probe-based approach". The "Hybrid rate control" unites a few of the most excellent characteristics of "receiver-based" and "transmitter-based control" in conditions of service elasticity and bandwidth effectiveness [39], [40]. Although will be able to merely go after the "Probe-based approach". One benefit of the "model-based approach" to "Probe-based approaches" for video multicasting is that it does not necessitate the

substitute of information between the groups beneath the "Probe-based approach". As a result, it removes the processing of every recipient and the bandwidth convention allied through the substitute of information.

- 2) Error control: It acquires the appearance of "Forward Error Correction (FEC)", "Delay Limit Retransmission", "Error Recovery or Error Hiding". There are three FECs: "channel coding", "source coding based FEC", and "joint source/channel coding". The benefit of all FEC methods above TCP is the decline of video communication delay. FEC based source coding be able to accomplish minor latency than "channel coding", and "joint source/channel coding" can accomplish the best presentation in the sense of rate deformation. The inconveniences of every one FEC methods are increased communication rate, different loss of flexibility characteristics. The "Feedback approaches" be able to utilized to increase FEC flexibility. Dissimilar FEC, which includes reiterating to make progress commencing a fail that may not happen, the "retransmission-based scheme" simply retransmit the failed packet. Therefore, "retransmission-based schemes" adapt to different failed features, ensuing in competent utilize of network sources.

This review concludes that future research efforts based on TCP congestion control limits [3], [5], [20] that have been identified can be developed with new enhancements to support end-to-end classical congestion control that will improve multimedia streaming End-to-end latency minimizes network overhead. At the same time, the problem of discarding tail queues also needs to be solved by a different routing-assisted congestion control algorithm and the new TCP fair queuing mechanism. The fair queuing mechanism [12] equitably distributes the bandwidth of the line to provide a relatively small queue delay for short-time communications to increase throughput.

VI. CONCLUSION

The Internet streaming is nowadays an important application utilized by general internet users. Nevertheless, the most excellent-attempt network is differentiated through active and impulsive changes in accessible bandwidth, which unfavorably affects video superiority. Therefore, it is significant that real-time recognition approaches with varying bandwidths ensure that the video adapts to the available bandwidth and transmits at the highest quality. The traditional view is to rely on end-user applications to deploy CCA to accomplish high network exploitation and a certain extent of traffic equality. In this article, we will discuss a systematic and comprehensive overview of TCP congestion control for TCP live streaming, as well as a study of the Internet streaming media mechanism. It ensures that the impact of congestion can seriously affect the fairness of the danger, and even data routing collapse. The router-based queue management scheme can effectively promote fairness goals and manage network congestion to share network resources fairly.

REFERENCES

- [1] G. Carlucci, L. De Cicco, S. Holmer, S. Mascolo, "Congestion Control for Web Real-Time Communication", *IEEE/ACM Transactions on Networking* Vol. 25, PP. 2629 - 2642, 2017.
- [2] S. D'Aronco, L. Toni, S. Mena, X. Zhu, P. Frossard, "Improved Utility-Based Congestion Control for Delay-Constrained Communication", *IEEE/ACM Transactions on Networking*, Vol. 25, PP. 349 - 362, 2017.
- [3] J. Luo, J. Jin, Feng Shan, "Standardization of Low-Latency TCP with Explicit Congestion Notification: A Survey", *IEEE Internet Computing*, Vol. 21, PP. 48 - 55, 2017.
- [4] R. Lübben, M. Fidler, "Service Curve Estimation-Based Characterization and Evaluation of Closed-Loop Flow Control", *IEEE Transactions on Network and Service Management*, Vol. 14, PP. 161 - 175, 2017.
- [5] Y. G Zhao, B. Zhang, C. Li, C. Chen, "ON/OFF Traffic Shaping on the Internet: Motivation, Challenges, and Solutions" *IEEE Network*, Vol. 31, PP. 48 - 57, 2017.
- [6] A. Biernacki, "Analysis of aggregated HTTP-based video traffic" *Journal of Communications and Networks*, Vol. 18, PP. 826 - 836, 2016.
- [7] L. De Cicco, Gaetano Carlucci, Saverio Mascolo "Congestion Control for WebRTC: Standardization Status and Open Issues", *IEEE Communications Standards Magazine*, Vol. 1, PP. 22 - 27, 2017.
- [8] Y. Li, H. Liu, W. Yang, D. Hu, X. Wang, Wei Xu, "Predicting Inter-Data-Center Network Traffic Using Elephant Flow and Sublink Information", *IEEE Transactions on Network and Service Management*, Vol. 13, PP. 782 - 792, 2016.
- [9] G. Tian, Y. Liu, "Towards Agile and Smooth Video Adaptation in HTTP Adaptive Streaming", *IEEE/ACM Transactions on Networking*, Vol. 24, PP. 2386 - 2399, 2016.
- [10] B. Briscoe, A. Brunstrom, A. Petlund, D. Hayes, et. al., "Reducing Internet Latency: A Survey of Techniques and Their Merits", *IEEE Communications Surveys & Tutorials*, Vol. 18, PP. 2149 - 2196, 2016.
- [11] K. Bilal, A. Erbad, "Edge computing for interactive media and video streaming", *Second International Conference on Fog and Mobile Edge Computing (FMEC)*, PP. 68 - 73, 2017.
- [12] G. Abbas, Z. Halim, Z. Haq Abbas, "Fairness-Driven Queue Management: A Survey and Taxonomy", *IEEE Communications Surveys & Tutorials*, Vol. 18, PP. 324 - 367, 2016.
- [13] A. Javadtalab, M. Semsarzadeh, A. Khanchi, S. Shirmohammadi, A. Yassine, "Continuous One-Way Detection of Available Bandwidth Changes for Video Streaming Over Best-Effort Networks", *IEEE Transactions on Instrumentation and Measurement*, Vol. 64, PP. 190 - 203, 2015.
- [14] Q. M. Qadir, A. A. Kist, Z. Zhang, "A Novel Traffic Rate Measurement Algorithm for Quality of Experience-Aware Video Admission Control", *IEEE Transactions on Multimedia* Vol. 17, PP. 711 - 722, 2015.
- [15] D. Li, Mingwei Xu, Ying Liu, Xia Xie, Yong Cui, Jingyi Wang, Guihai Chen "Reliable Multicast in Data Center Networks", *IEEE Transactions on Computers*, Vol. 63, PP. 2011 - 2024, 2014.
- [16] Z. Li, Xiaoqing Zhu, Joshua Gahm, Rong Pan, Hao Hu, Ali C. Begen, David Oran "Probe and Adapt: Rate Adaptation for HTTP Video Streaming At Scale", *IEEE Journal on Selected Areas in Communications*, Vol. 32, PP. 719 - 733, 2014.
- [17] S. A. Memon, S. R. Hassan, N. A. Memon, "Evaluation of video streaming performance over the peer-to-peer network", *International Conference on Collaboration Technologies and Systems (CTS)*, PP. 413 - 420, 2014.
- [18] E. Grigorescu, C. Kulatunga, G. Fairhurst, "Evaluation of the impact of packet drops due to AQM over capacity limited paths", *21st IEEE International Conference on Network Protocols (ICNP)* PP. 1 - 6, 2013.
- [19] M. Gorius, Y. Shuai, T. Herfet, "Dynamic media streaming under predictable reliability", *IEEE International Symposium on Broadband Multimedia Systems and Broadcasting*, PP. 1 - 6, 2012.
- [20] P. Yang, L. Xu, "A survey of deployment information of delay-based TCP congestion avoidance algorithm for transmitting multimedia data", *IEEE GLOBECOM Workshops (GC Wkshps)* PP. 18 - 23, 2011.
- [21] C. Perkins and V. Singh, "Multimedia Congestion Control: Circuit Breakers for Unicast RTP Sessions", *RFC 8083*, RFC Editor, Mar. 2017.
- [22] B. Briscoe et al., "Reducing Internet latency: A survey of techniques and their merits", *IEEE Communications Surveys Tuts.*, vol. 18, no. 3, pp. 2149-2196, 3rd Quart., 2016.
- [23] Y. Zaki, T. Pötsch, J. Chen, L. Subramanian, and C. Görg, "Adaptive congestion control for unpredictable cellular networks", in *Proc. ACM SIGCOMM*, vol. 45. no. 5, pp. 509-522, 2015.
- [24] R. Pan et al., "Pie: A lightweight control scheme to address the bufferbloat problem", in *Proc. IEEE HPSR*, pp. 148-155, 2013.
- [25] J. Jiang, V. Sekar, and H. Zhang, "Improving Fairness, Efficiency, and Stability in HTTP-Based Adaptive Video Streaming with Festive", *Proc. ACM CoNEXT '12*, pp. 97-108, 2012.
- [26] K. Nichols and V. Jacobson, "Controlling queue delay", *Queue ACM*, vol. 10, no. 5, pp. 20:20-20:34, May 2012.
- [27] K. Jacobsson, L. L. H. Andrew, A. Tang, S. H. Low, and H. Hjalmarrsson, "An improved link model for window flow control and its application to FAST TCP", *IEEE Trans. Autom. Control*, vol. 54, no. 3, pp. 551-564, Mar. 2009.
- [28] D. A. Hayes and G. Armitage, "Revisiting TCP congestion control using delay gradients", in *Proc. 10th Int. Conf. Res. Netw.*, vol. 2, pp. 328-341, 2011.
- [29] L. De Cicco, G. Carlucci, and S. Mascolo, "Understanding the dynamic behavior of the Google congestion control for RTCWeb", In *Proc. Int. Packet Video Workshop (PV)*, pp. 1-8, 2013.
- [30] A. Mansy, B. V. Steeg, and M. Ammar, "Sabre: A Client-Based Technique for Mitigating the Buffer Bloat Effect of Adaptive Video Flows", *Proc. ACM MMSys '13*, pp. 214-25, 2013.
- [31] D. Wischik and N. McKeown, "Part I: Buffer sizes for core routers", *ACM SIGCOMM Comput. Communications Rev.*, vol. 35, no. 3, pp. 75-78, Jul. 2005.
- [32] G. Raina, D. Towsley, and D. Wischik, "Part II: Control theory for buffer sizing", *ACM SIGCOMM Comput. Communications Rev.*, vol. 35, no. 3, pp. 79-82, Jul. 2005.
- [33] M. Enachescu, Y. Ganjali, A. Goel, N. McKeown, and T. Roughgarden, "Part III: Routers with very small buffers", *ACM SIGCOMM Comput. Communications Rev.*, vol. 35, no. 3, pp. 83-90, Jul. 2005.
- [34] K. Ramakrishnan, S. Floyd, and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP, document RFC 3168, Internet Requests for Comments, Sep. 2001.
- [35] R. Jain, "A delay-based approach for congestion avoidance in interconnected heterogeneous computer networks", *ACM SIGCOMM Comput. Communications Rev.*, vol. 19, no. 5, pp. 56-71, Oct. 1989.
- [36] L. S. Brakmo and L. L. Peterson, "TCP Vegas: End to end congestion avoidance on a global Internet", *IEEE J. Sel. Areas Communications*, vol. 13, no. 8, pp. 1465-1480, Oct. 1995.
- [37] E. Brosh, S. A. Baset, V. Misra, D. Rubenstein, and H. Schulzrinne, "The delay-friendliness of TCP for real-time traffic", *IEEE/ACM Trans. Netw.*, vol. 18, no. 5, pp. 1478-1491, 2010.
- [38] Z. Li et al., "Probe and Adapt Rate Adaptation for HTTP Video Streaming at Scale", *IEEE JSAC*, vol. 32, pp. 719-33, 2014.
- [39] M. Ghobadi et al., "Trickle: Rate limiting YouTube Video Streaming", *Proc. USENIX ATC '12*, pp. 191-96, 2012.
- [40] S. Akhshabiet al., "What Happens When HTTP Adaptive Streaming Players Compete for Bandwidth?", *Proc. ACM NOSSDAV '12*, pp. 9-14, 2012.